

The Thermal Wall:

Where it came from and how to live with it?

Ronny Ronen
Senior Principal Engineer
Director of Architecture Research
Microprocessor Technology Labs
Haifa, Israel

Intel Corporation

**The 10th Intel
EMEA Academic forum
Gdansk, Poland
May 19, 2005**

Contributors: Lev Finkelstein, Eli Savransky, Efi Rotem



Point Mugu (USA) Airshow 2004

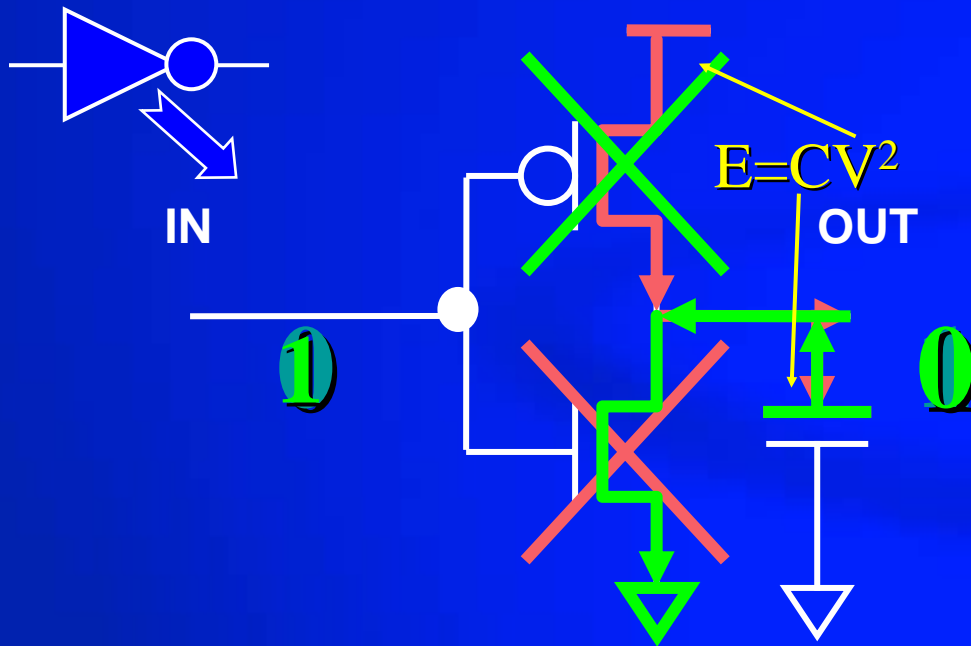
[http://www.richard-seaman.com/Aircraft/
AirShows/PointMugu2004/Highlights/](http://www.richard-seaman.com/Aircraft/AirShows/PointMugu2004/Highlights/)

Agenda

- **Background**
- **Process Scaling**
- **The Thermal Challenge**
- **Variations**
- **Dynamic Thermal Management**
- **Conclusions**

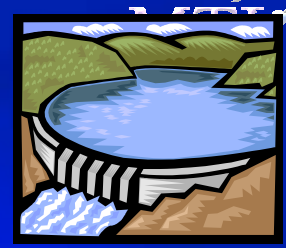
Power and the digital world (1)...

- Power is consumed:
 - When capacitance is charged and discharged.
 - A charged cap is a logical '1', a discharged cap is '0'.

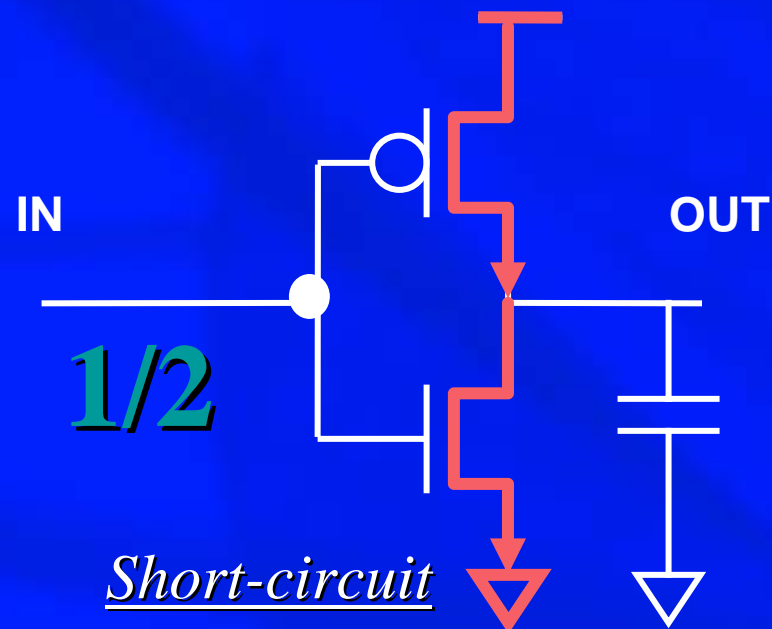
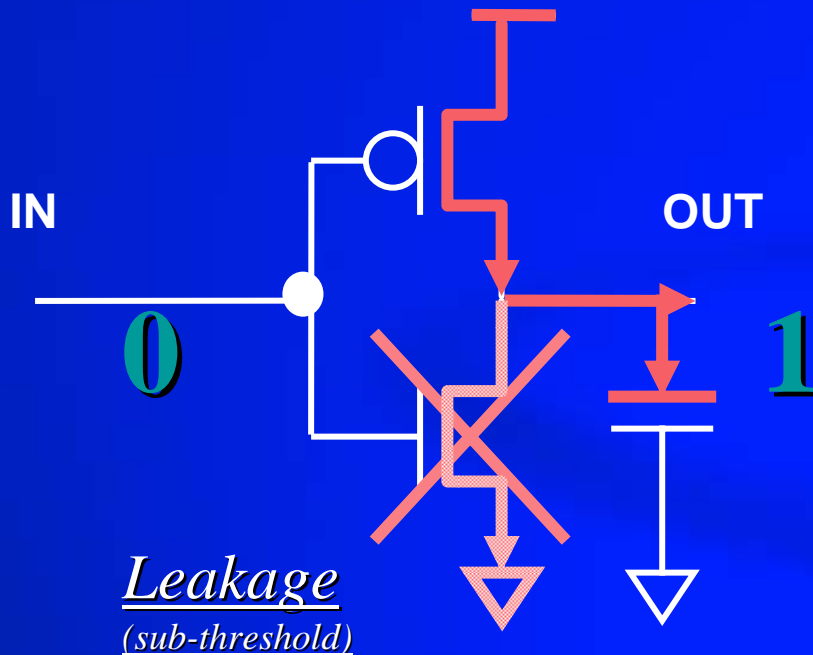


- Power is also consumed on wires – busses, interconnects.

Power and the digital world (2)...

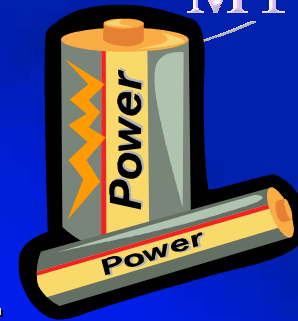


- Secondary effects like leakage and short-circuit current are increasing with advanced process technologies.



- Leakage became a 1st order factor
 - Used to be 7-15% not long ago, 20%-50% in coming generations.

Power & Energy



Energy

- Relevant in respect to:
 - Battery life - lower energy per task → longer battery life.
 - Electric bills - lower energy per task → lower bills.
- $E = \alpha CV^2$ - Proportional to capacitance, voltage², and activity

Power

- *Energy per time* - measured in Watts.
- $P = \alpha CV^2f$ (α : activity, C: capacitance, V: voltage, f: frequency)
- Peak power is a concern:
 - Higher power → higher current.
 - Higher power → higher temperature.

Power Density

- Think of watts/cm².
- Higher power, Smaller area → Higher power density
- Denser power → Harder to cool

Thermal



Thermal

- Power → heat → higher temperature
- Temperature increased with
 - Higher power
 - Higher power density
 - Higher ambient temperature
 - Weaker cooling
(cooling = form factor, cost)
- Higher temperature risks
 - Immediate malfunctioning
 - Reduced long term reliability
 - Higher leakage power
 - Noise





Cause and Effect

- More transistors, bigger transistors, higher voltage, higher activity
→ Higher energy, higher power
- Higher frequency
→ Higher power
- Denser logic, higher power, smaller die
→ Higher power density
- Higher power, higher power density
→ Higher temperature
- Higher temperature
→ Higher leakage power
→ Higher temperature...

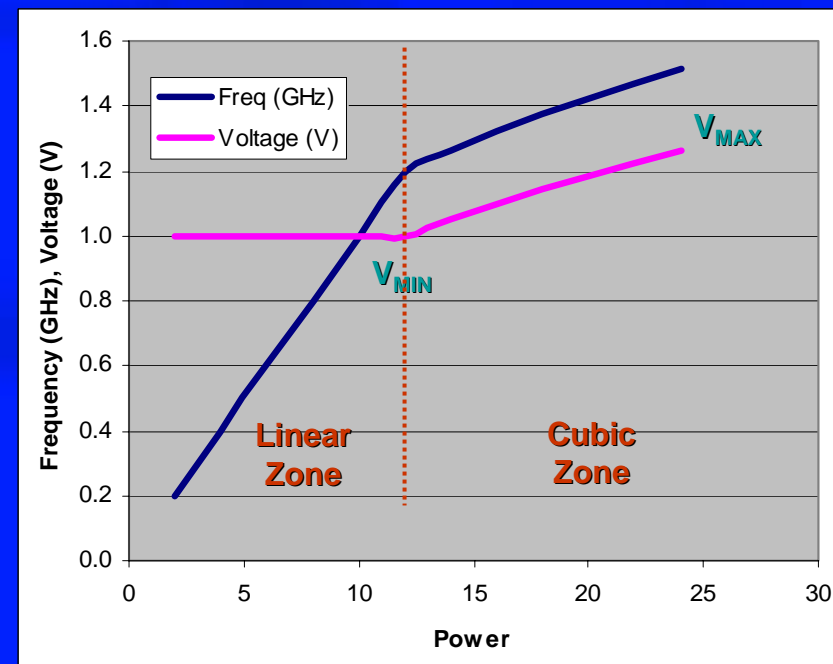
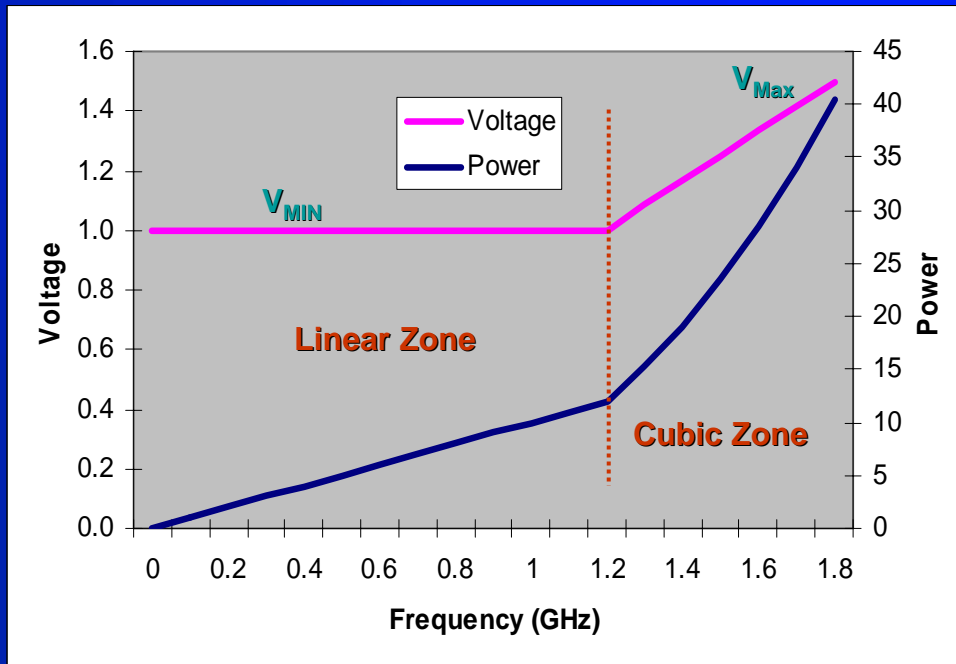
Performance, Frequency, Transistors

- **Performance = IPC × Frequency**
- **Sources of frequency increase**
 - Process: lower delay time
 - Microarchitecture: deeper pipeline
- **Higher frequency → Lower IPC**
 - Memory access do not scale ...
 - e.g., 50% frequency gain → 35%-45% performance
- **More transistors → Higher IPC. Typical#**
 - Single core: 2X transistors → 1.4 IPC
 - Multi core: 2X transistor → 1.8 IPC

Performance, Power

Physics: Voltage/Frequency scaling

- $Freq \approx k * Voltage$ ($F \approx kV$)
 - Within a limited voltage range V_{MIN}/V_{MAX}
- **Active Power = Activity * Capacitance * Voltage² * Freq** ($P = \alpha CV^2 f$)
 - ➔ Power is proportional to f^3 ($>V_{MIN}$)
 - ➔ Power is proportional to f ($\leq V_{MIN}$)
- **Used for Static and Dynamic Voltage/Frequency) Scaling (DVS)**

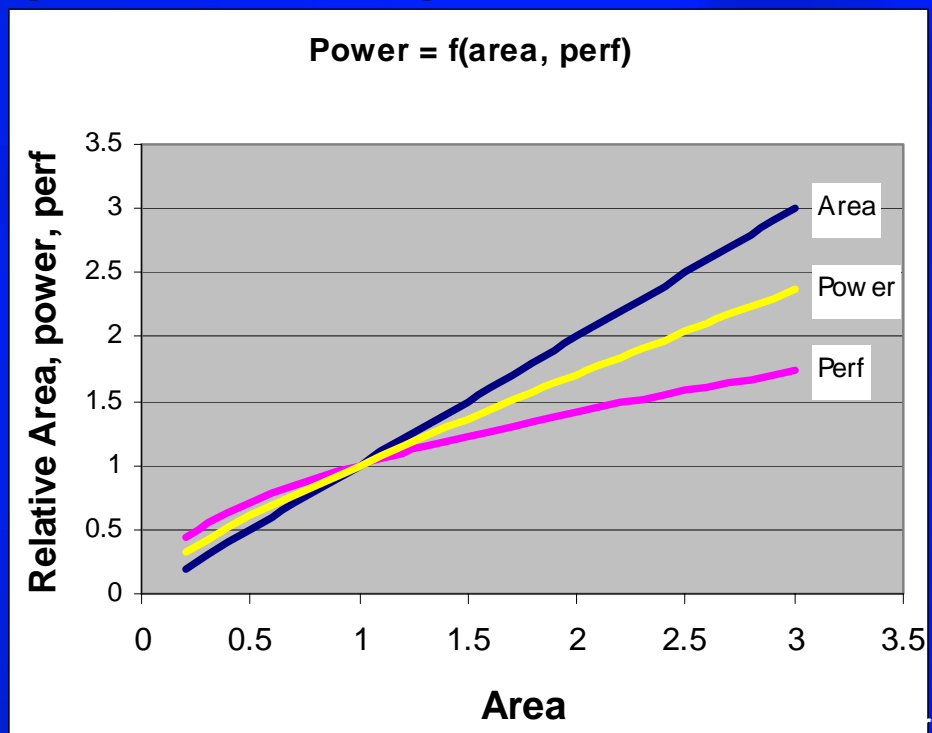


Performance, Power, Area

Empirics – the power growth

$Power = f(Perf, Area)$

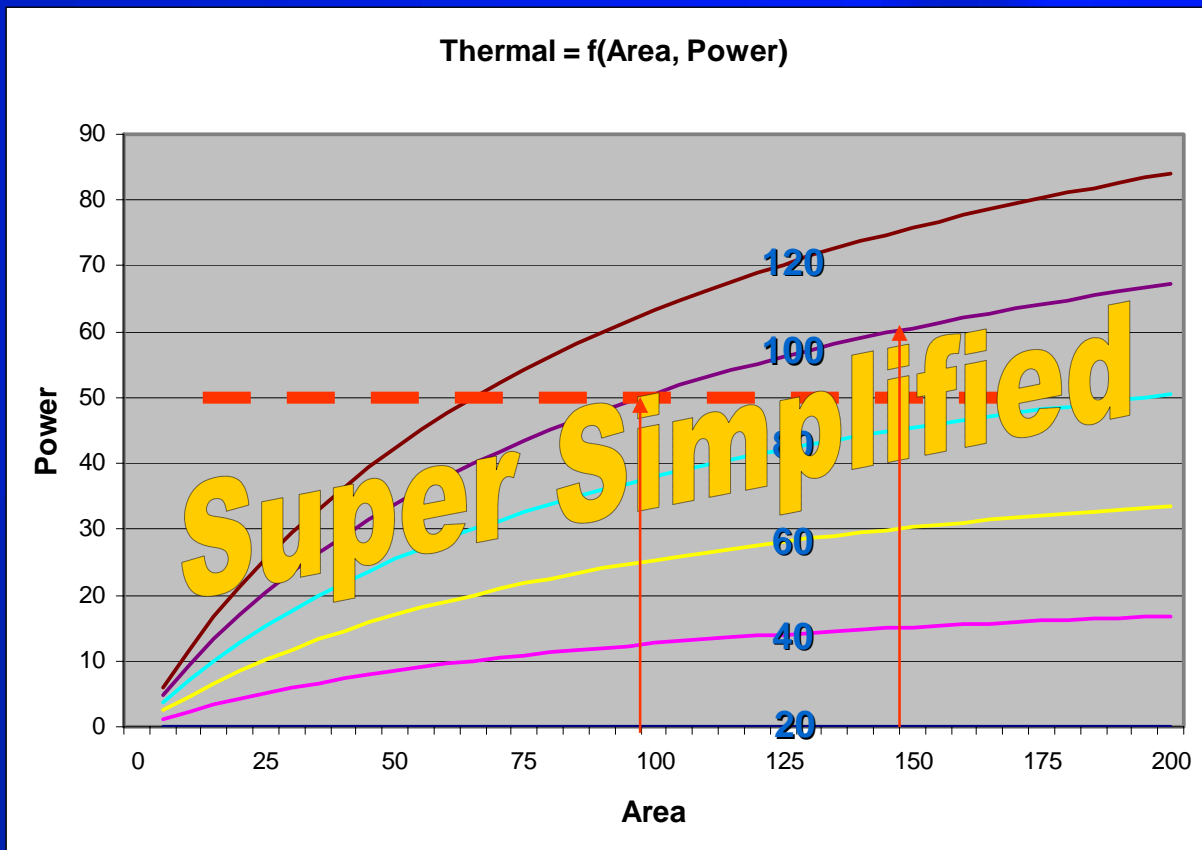
- $Capacitance = k * Area$ ($C=kA$)
- Power Growth Guesstimate (same process, same voltage)
 - Area growth \rightarrow power growth
 - Performance growth \rightarrow power growth
 - ➔ Empirically, for the same process and voltage, power growth is ~ average of both (combined effect of αC)



Power, Power Density, Thermal

Thermal = f(Power, Power Density) = f(Power, Area)

- Higher Power → Higher Temperature
- Smaller Area → Higher Temperature



Note:
Same power
smaller area
→ Higher temp.

Qualitative Data
Assumes ambient temp. 20°

Process Technology – In Practice

Every new process technology (~2-3 years cycle) changes the physical attributes of the transistor:

- Size reduced to 0.5X
- Delay time reduced to **0.75X**
- Switch energy reduced by **1.75X** (Relative to CV^2)
 - Capacitance per transistor reduced to 0.7X
 - Operation voltage reduced to **0.9X**

Two extreme scenarios:

- Ideal “Shrink”
 - Same μ arch

1X	#Xistors
0.5X	size (0.7X per dimension)
1.35X	frequency
0.9X ¹	IPC
0.75X	peak power
1.2X	performance (IPC×freq)
1.5X	power density

- Ideal New design
 - Same die size

2X	#Xistors
1X	size
1.35X	frequency
1.3X-1.7X ¹	IPC
1.5X	peak power
1.8-2.3X	performance
1.5X	power density

¹ IPC Decreased in higher frequencies

In Pictures...

Process Technology

Processor 1.5 μ 1.0 μ 0.8 μ 0.6 μ 0.35 μ 0.25 μ 0.18 μ 0.13 μ 90nm 65nm

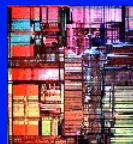
Intel386™ DX
Processor



Intel486™ DX
Processor



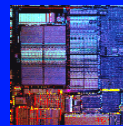
Pentium®
Processor



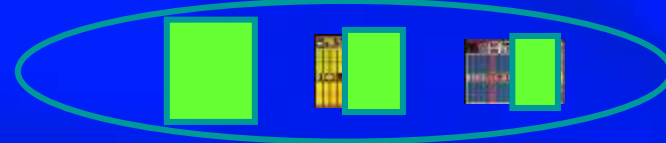
Pentium® Pro
Processor



Pentium® II
Processor



Pentium® III
Processor



Pentium® 4
Processor



Pentium® M
Processor

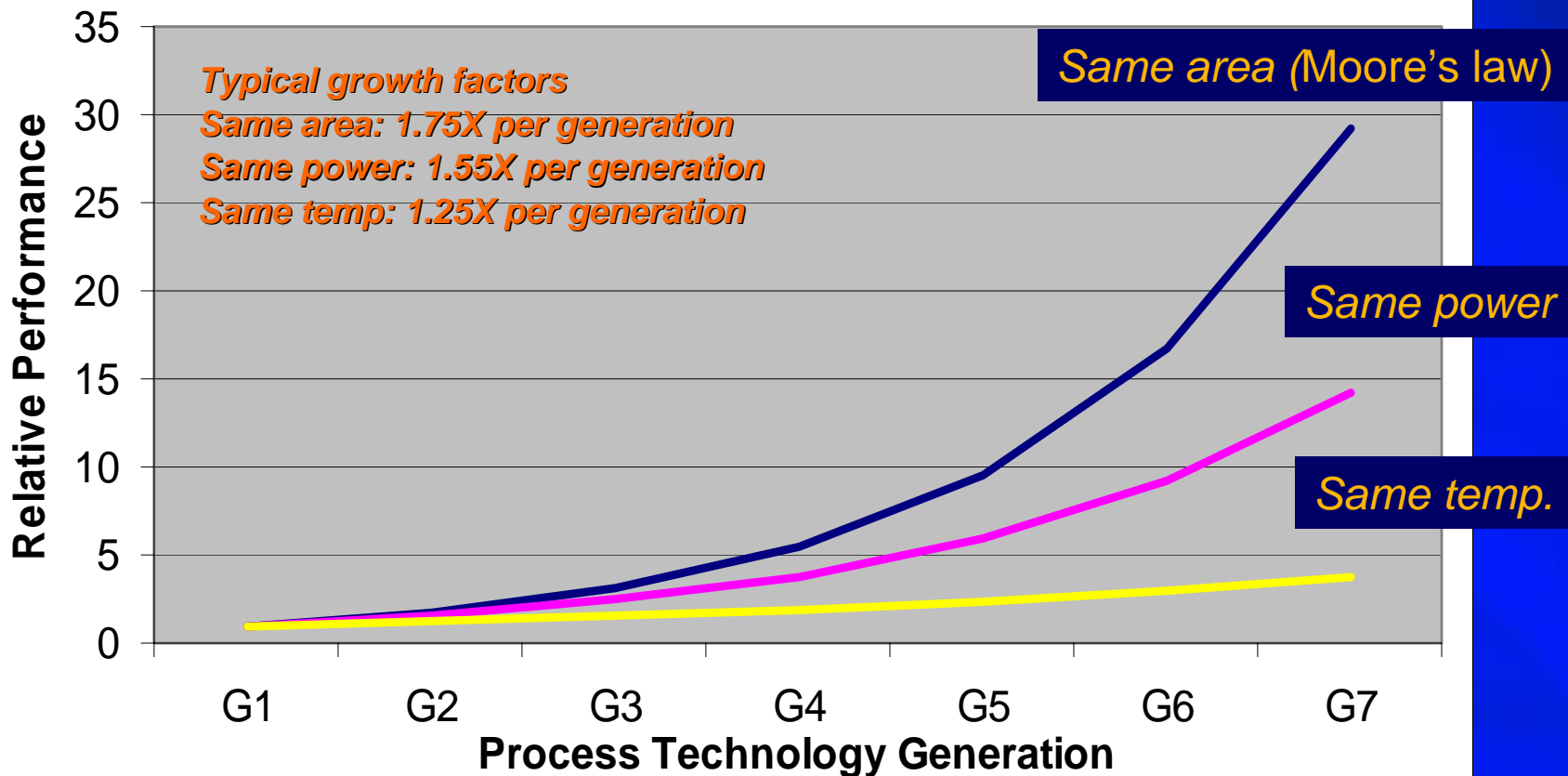


Next generation...



The Thermal Wall

Performance scaling of microprocessors



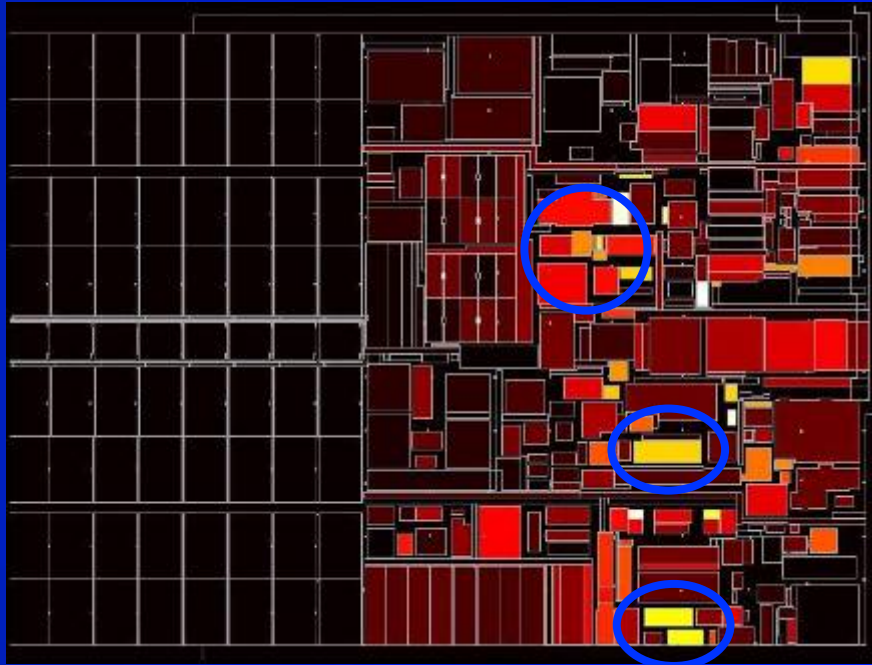
In a thermally limited environment

Evolutionary Uarch will diminish its performance return

And it Gets Even Worse...

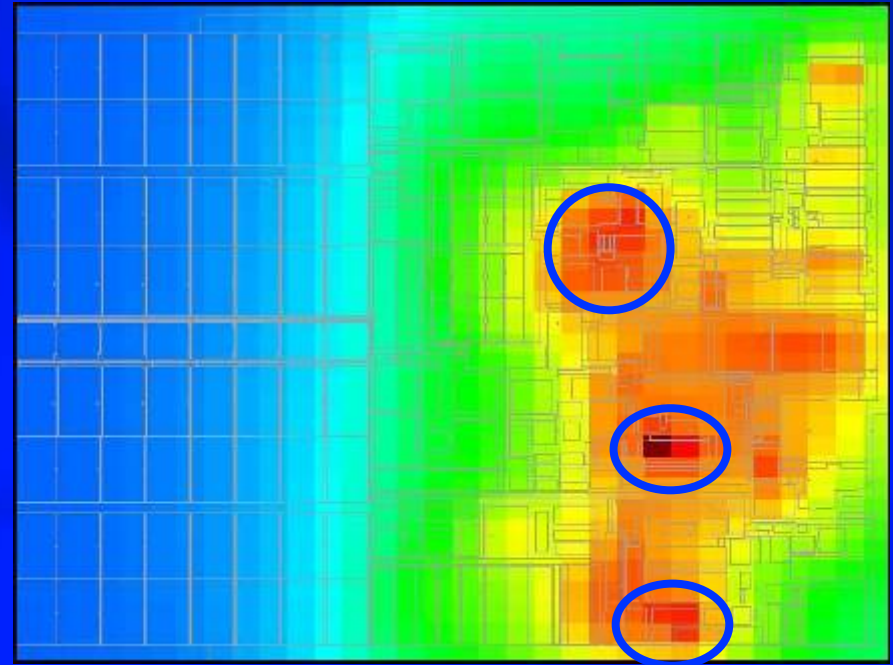
- **Variation in space**
- **Variation in time**
- **Variation in process**
- **Variation in platform (form factor)**
- **Variation in cooling solutions**
- **OS dependence (ACPI)**

Thermal Variation - Space



Power Density (simulated)¹

Color codes: (lowest) black, red, orange yellow, white (highest)



Thermal Map (simulated)²

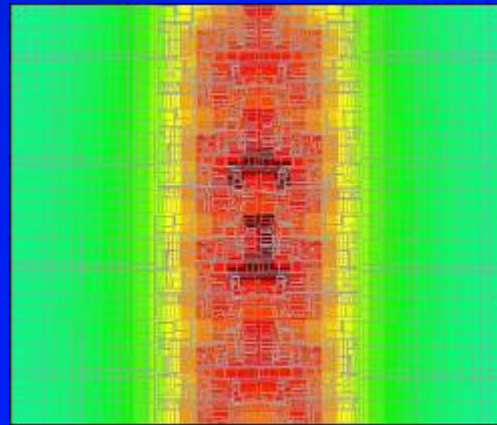
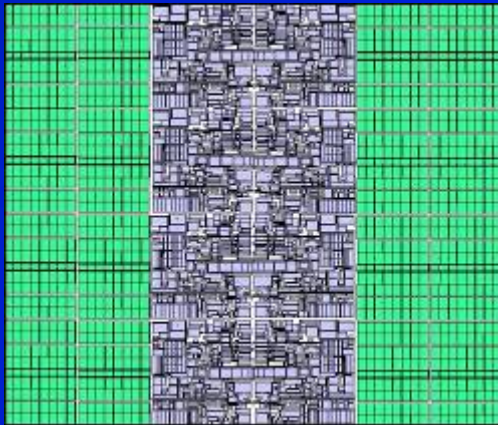
(lowest) blue, green, yellow, orange, red (highest)

Pentium M Processor power density example

¹ Source : Intel® Pentium® M Processor Power Estimation, Budgeting, Optimization, and Validation
Dani Genossar, Nachum Shamir, ITJ Q2/2003

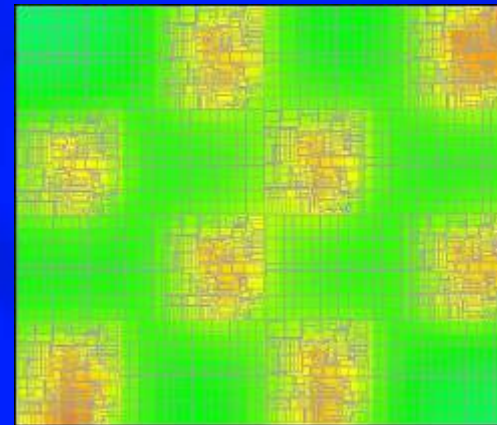
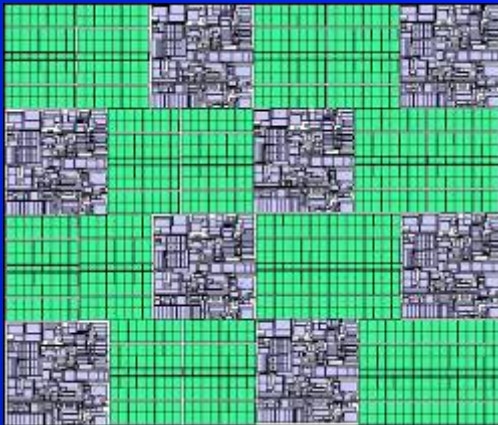
² Source: Lev Finkelstein, Intel 2005

Thermal Variation – Space / CMP



8 cores on the center
Tjmax: 100.4°C

Qualitative Data

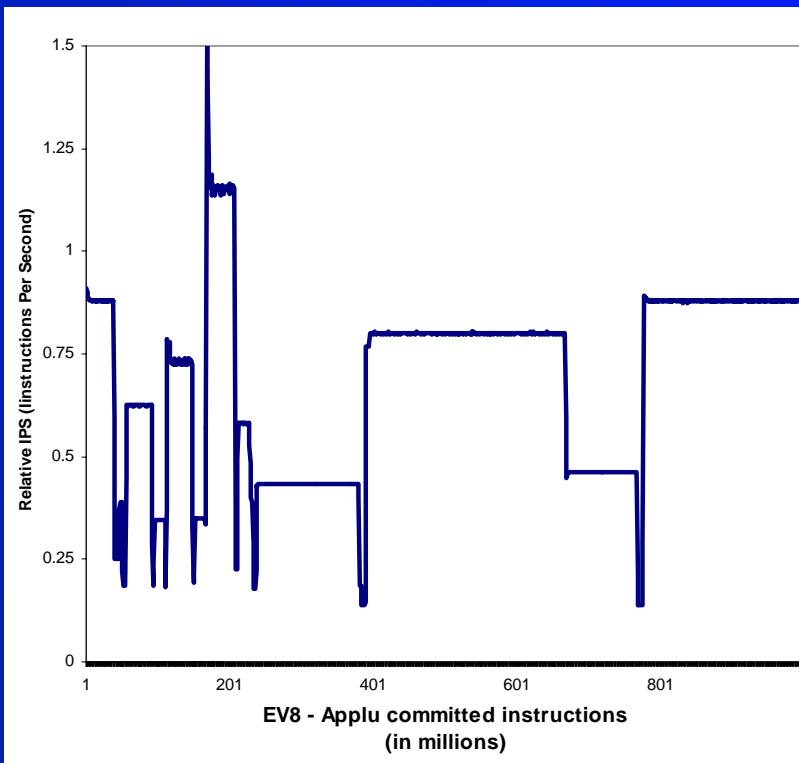


8 cores checkers layout
Tjmax: 96.3°C

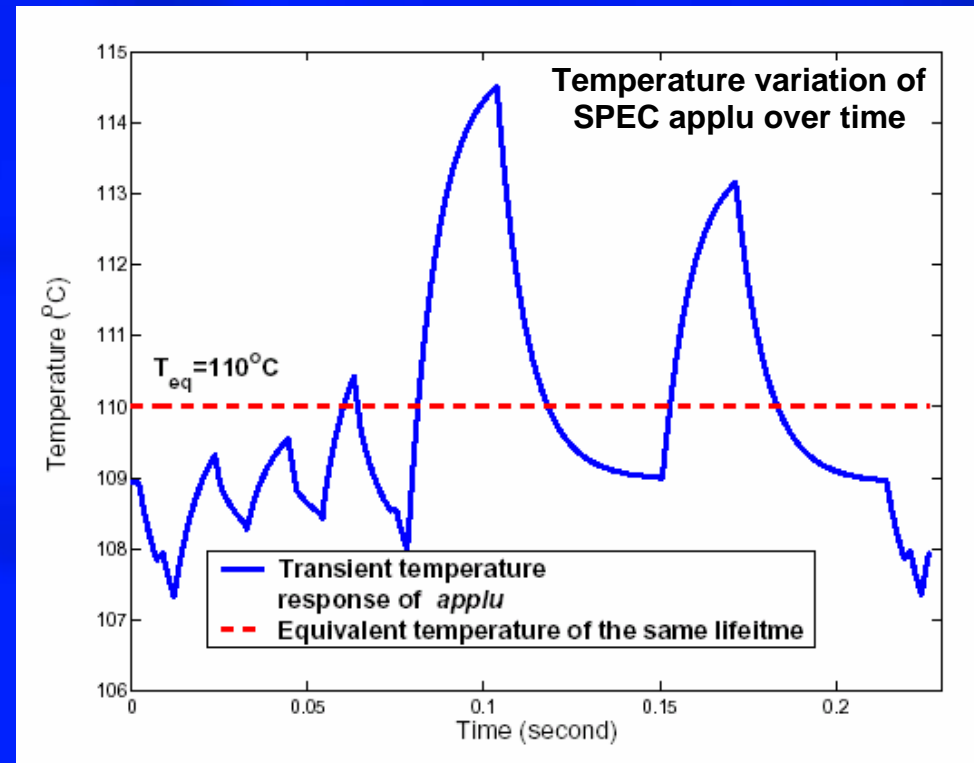
- Uniform power distribution is better
- Chip Multi-Processor (CMP) can exploit that

Thermal Variation – Time

- Follows activity → tracks IPC
- Varies among applications
- Varies within an application



Source: “Processor Power Reduction Via Single-ISA Heterogeneous Multi-Core Architectures”
Rakesh Kumar et.al Micro’03

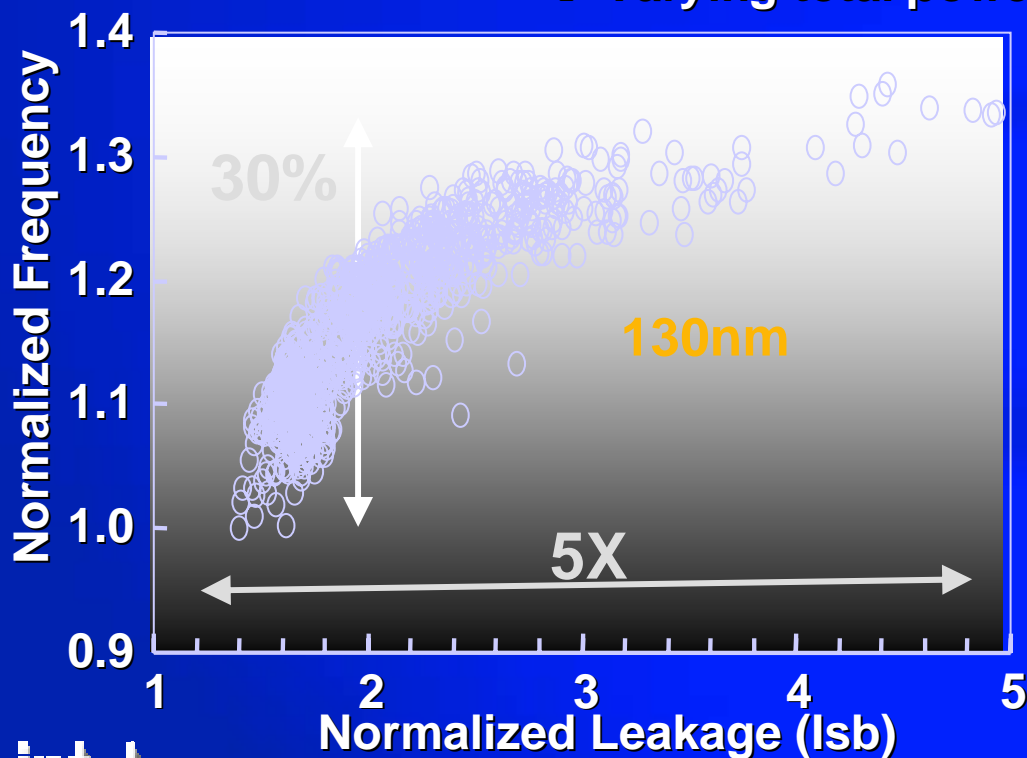


Source: “A Quick Thermal Tutorial”
Kevin Skadron, Mircea Stan, U. of Virginia 2005

Two graphs are not from exact same code

Thermal Variation - Process

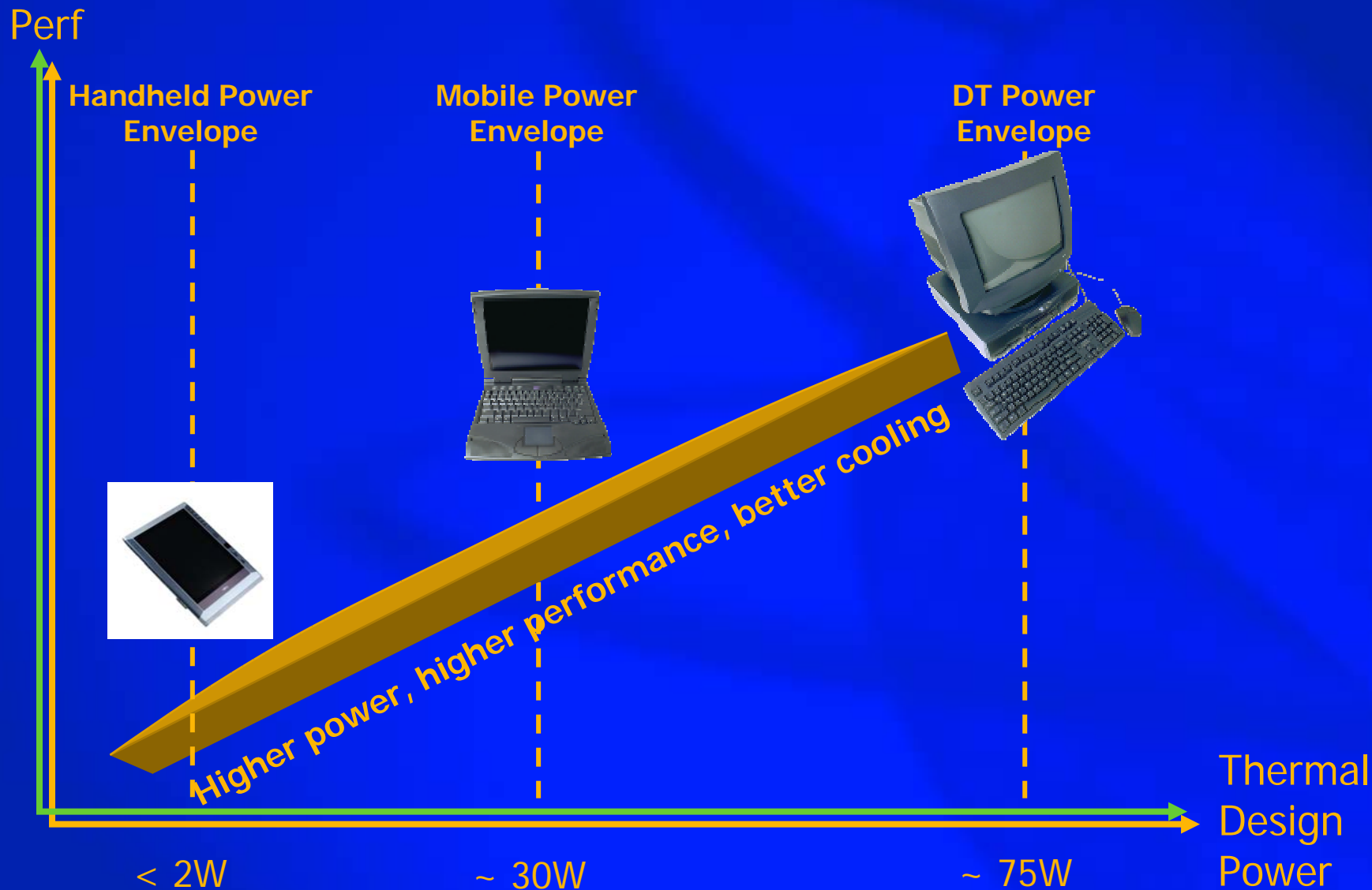
- Random dopant fluctuations
- Sub-wavelength Lithography
- ➔ Varying feature size
- ➔ Varying delay time
- ➔ Varying leakage power
- ➔ Varying total power, temperature



Frequency
~30%

Leakage Power
~5-10X

Thermal Variation - Platform



Implications

- Thermal events no longer a rare exception
 - Can happen frequently
 - In different scenarios
- Worst-case design is extremely inefficient and error prone

Approach:

- *Dynamic Thermal Management*

Dynamic Thermal Management (DTM)

- Basics
- “Linear” Thermal Throttling
- DVS* based Thermal Throttling
- Micro-architectural Thermal Throttling
- Advanced DTM
- Optimal DTM

*DVS: Dynamic Voltage Scaling

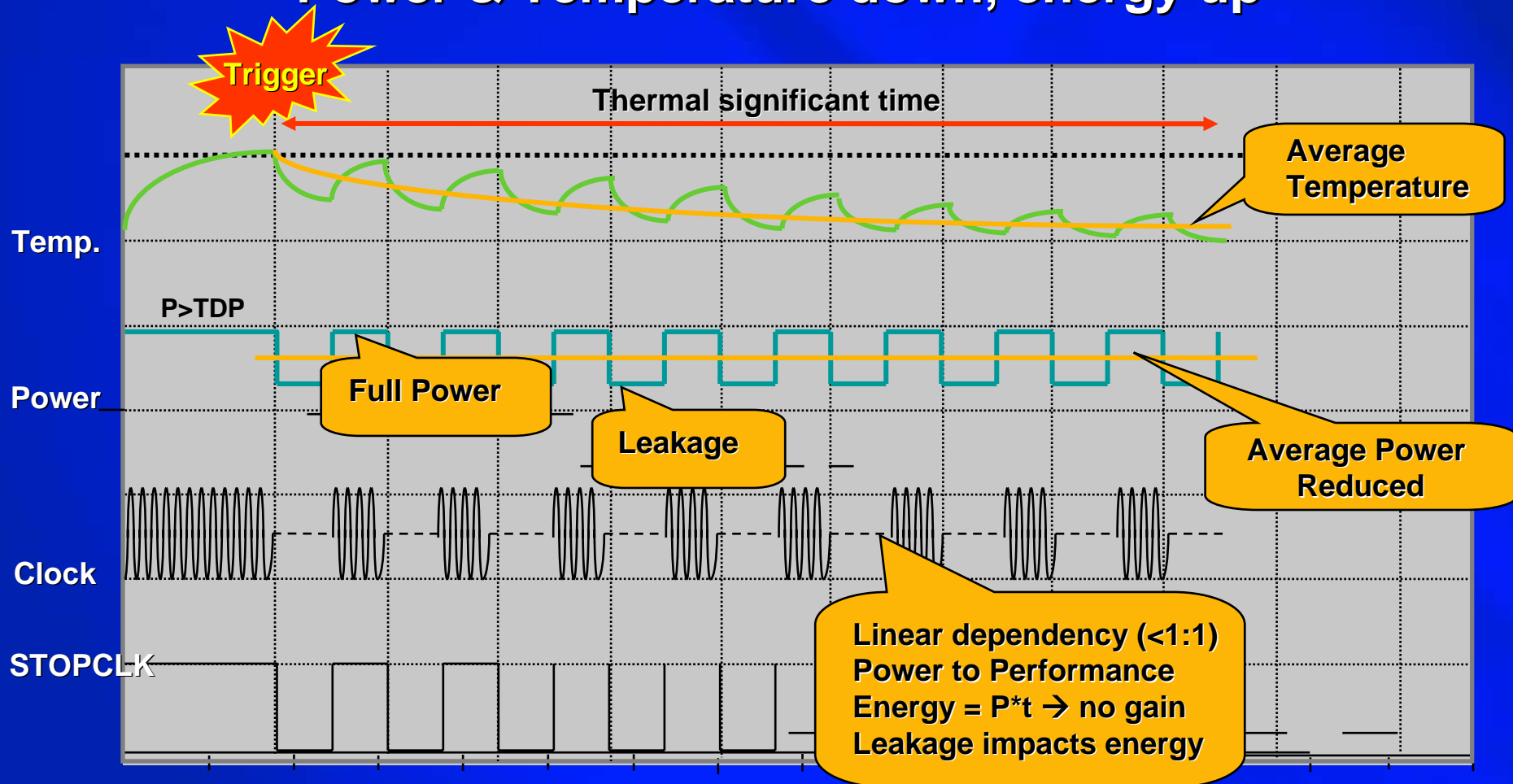
Thermal Management - Basics

- Avoid over heating by reduction in activity
- Triggered by thermal event
 - Actual or predicted
 - Traditionally detected by thermal sensor
- Activity reduction methods
 - “Stop clock”
 - Frequency reduction
 - Voltage/frequency scaling, ...
- When to throttle
- Where to throttle
- For how long? (when to stop)



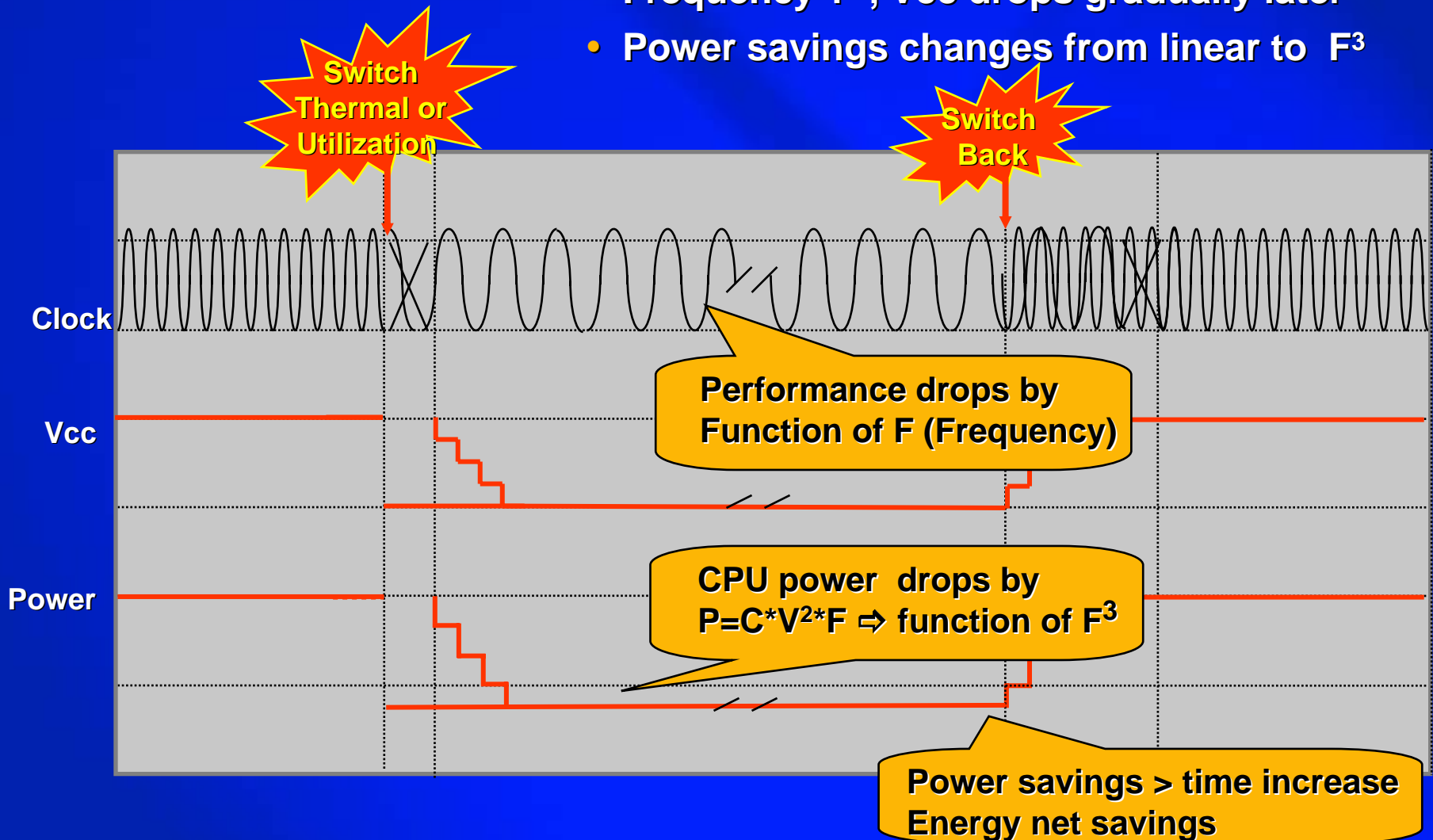
"Linear" Thermal Throttling

- Stop all activity for a predefined time
- Power & Temperature down, energy up



DVS based Thermal Throttling

- Reduce both Voltage & Frequency
- Frequency 1st, Vcc drops gradually later
- Power savings changes from linear to F^3

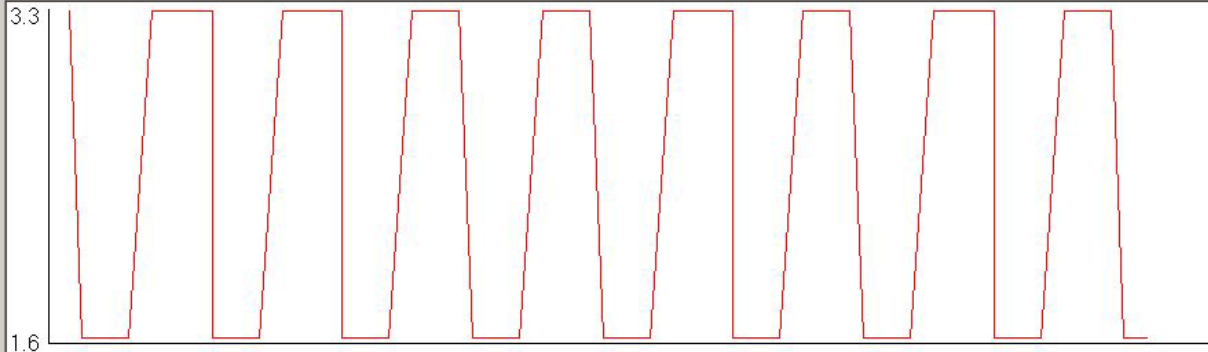


DVS vs. Linear Throttling

Throttling Demo

File

60



Frequency throttling

Temperature



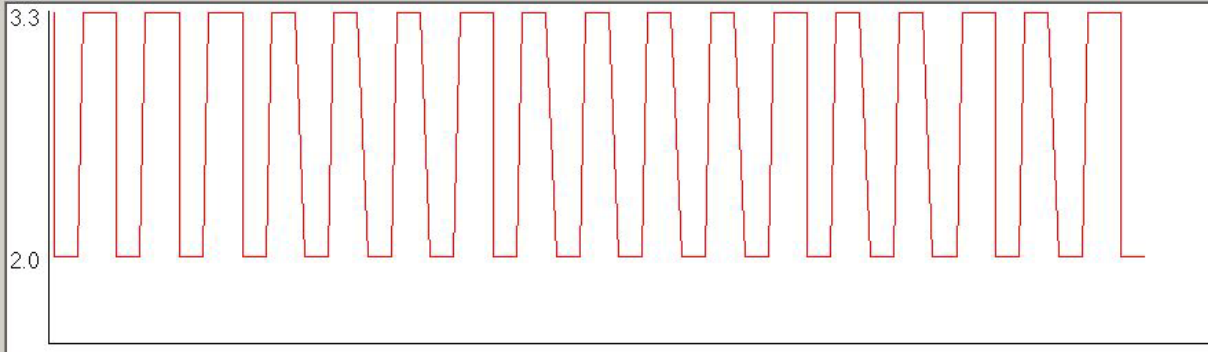
Frequency

2362

Performance (vs. fmax)



72%



V&f throttling

Temperature



Frequency

2701

Performance (vs. fmax)



82%

DVS Thermal Throttling – The Challenge

- Maximize performance for given thermal environment

Constraints

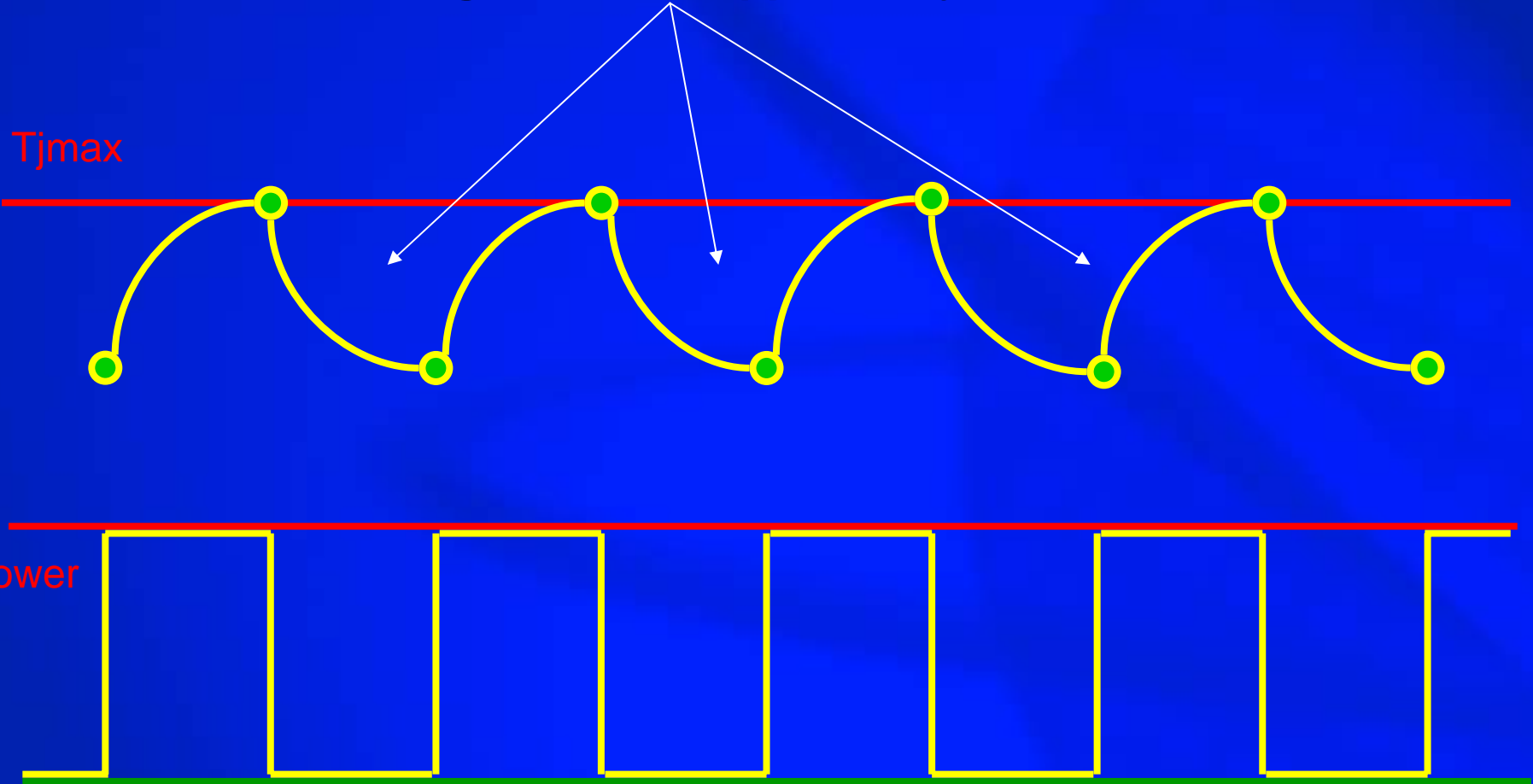
- Voltage/frequency changes take time
- Longer operation at lower frequency → lower performance

Questions

- When to throttle?
 - Should we wait for the last moment?
- Where to throttle?
 - Should we go for lowest voltage/frequency?
- When to stop throttling?
 - How much time should we cool down? *Sprint or Marathon?*
 - To which voltage/frequency level?

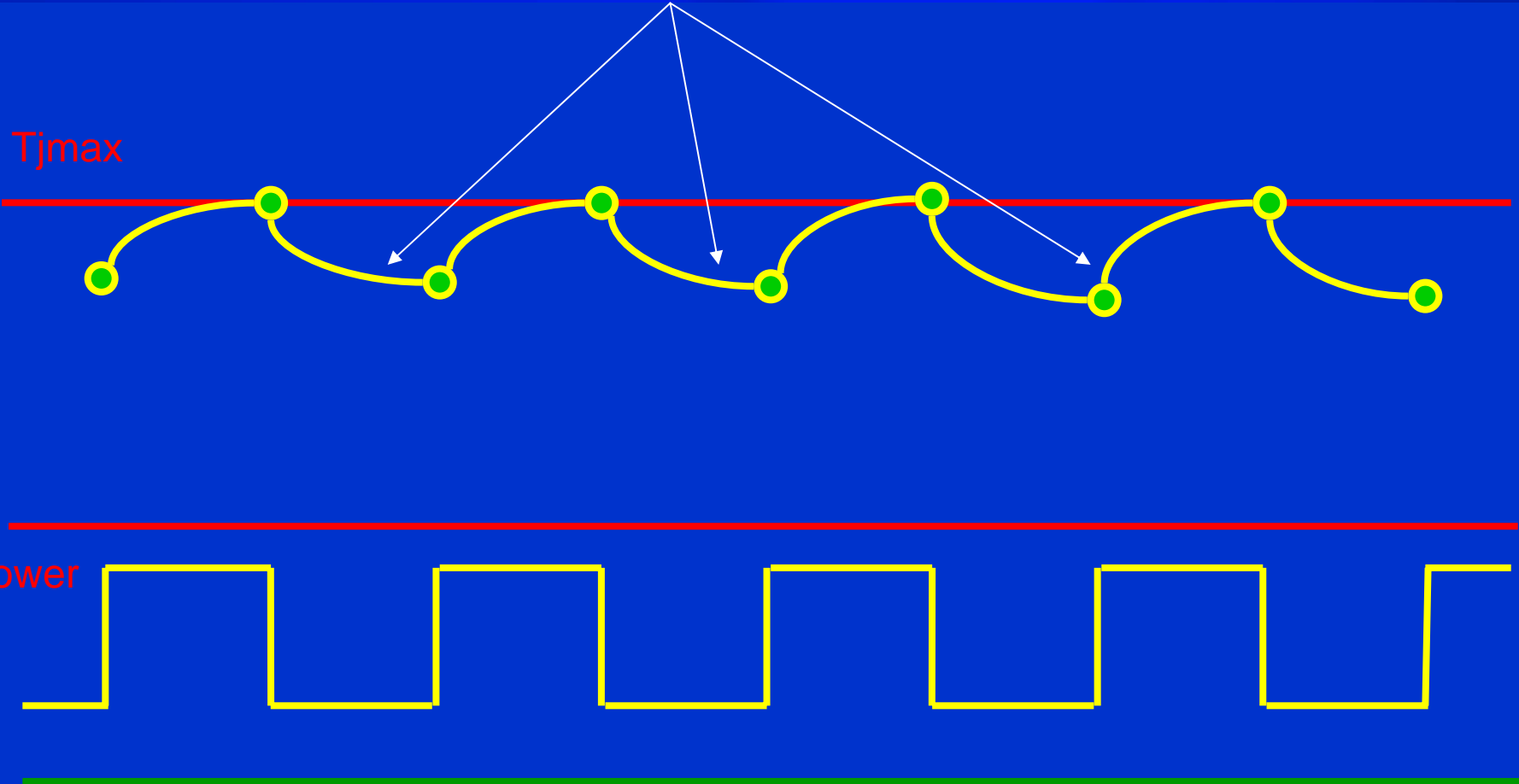
Where to Throttle?

Region of lost opportunity



Where to Throttle?

Somewhat Better



Where to Throttle?

Much Better

T_{jmax}



Power





Frequency throttling

Temperature



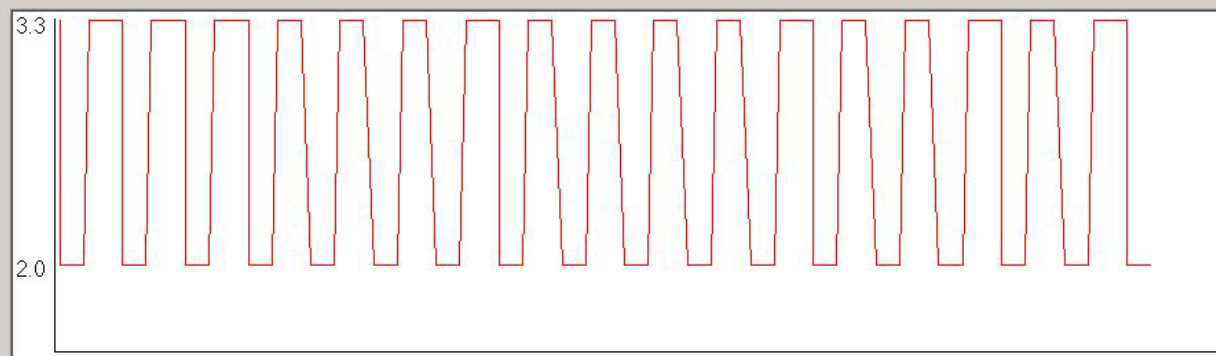
Frequency

2362

Performance (vs. fmax)



72%



V&f throttling

Temperature



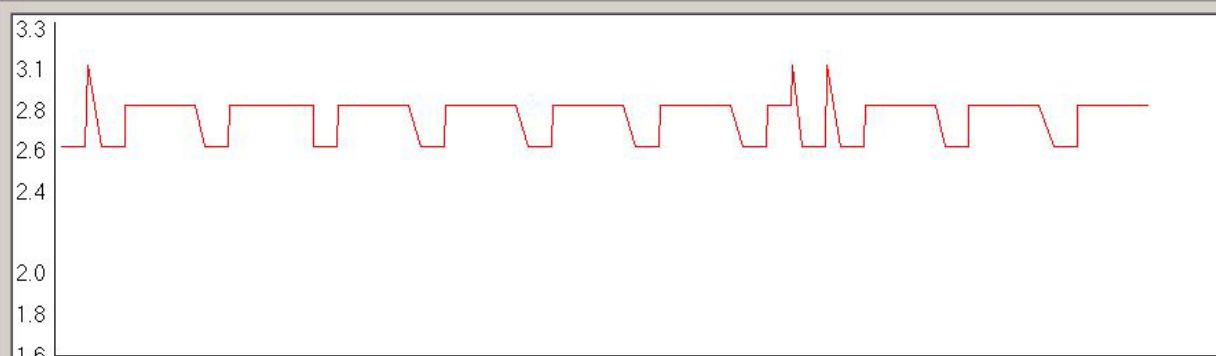
Frequency

2701

Performance (vs. fmax)



82%



Smart throttling

Temperature



Frequency

2789

Performance (vs. fmax)



85%

Micro-Architectural Throttling

- Reduce activity by slowing certain units
- Reduces power → reduces temperature
- Examples:
 - Stop or slow down instruction fetch¹
 - Stop speculation
 - Adaptively reduce array sizes
- Advantages
 - Fast activation
 - Super-linear savings
 - Focused targets
- Disadvantages
 - Usually, less efficient than DVS

¹ Pipeline Gating:

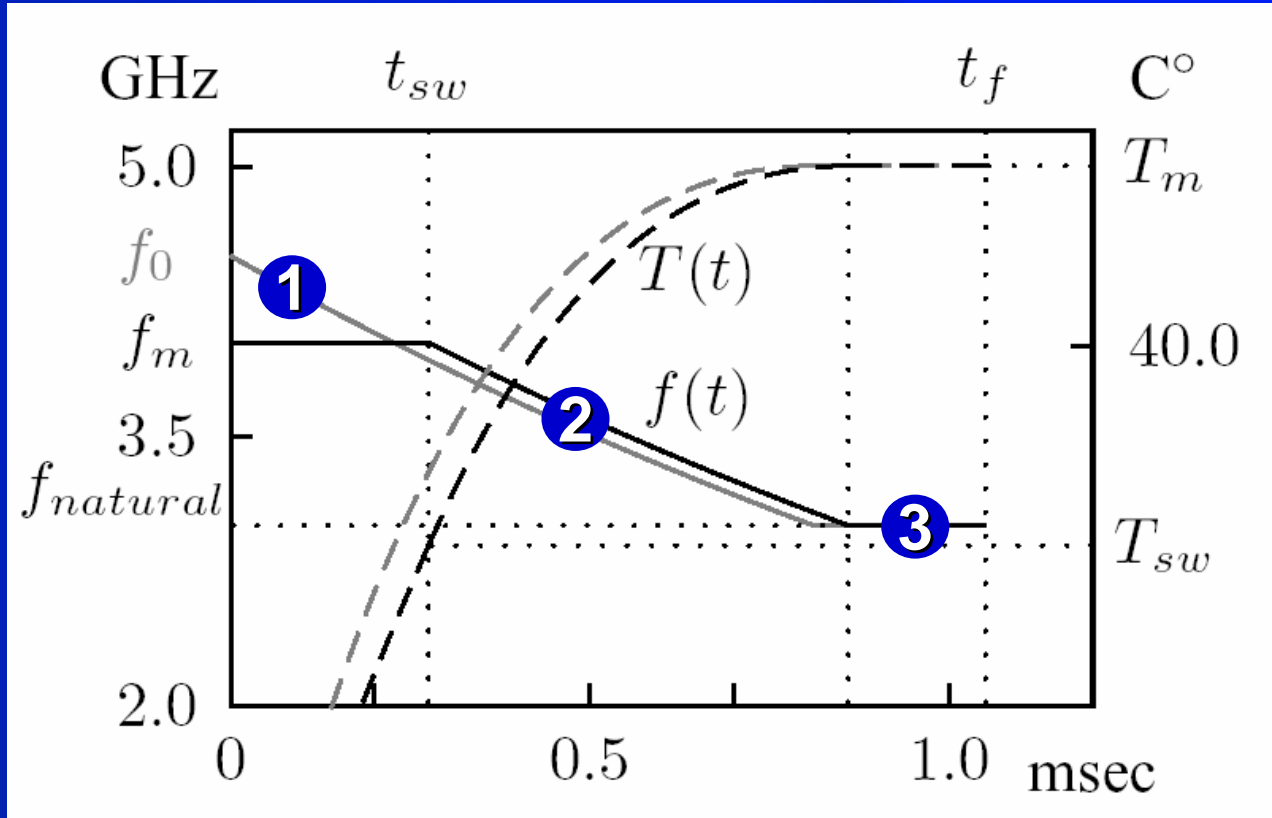
Speculation Control for Energy Reduction. Srilatha Manne, Artur Klauser, Dirk Grunwald, ISCA '98

Advanced DTM

- **Goal:**
Maximize performance per thermal constraints
- **Reactive systems: fast, minimal input,**
 - Throttling: on/off scheme, thermal interrupts.
- **Controllers**
 - PID¹ control - tune power to keep temperature
 - Fast, minimal history usage, needs power estimation
- **Even smarter mechanisms**
 - Predictive techniques
use history to predict power & temperature behavior
- **Tradeoff**
input type/size \Leftrightarrow computational efforts \Leftrightarrow
reaction time \Leftrightarrow quality

"Optimal" Control – example

Maximal Performance per Given Max T_j



1. Can run higher than maximal "natural" frequency and can dissipate more than maximal power for limited time
2. Start decreasing power even before reaching Max T_j
3. Eventually approach "natural" maximal power & frequency

Optimal strategy for a single-RC model¹
(Starting at an ambient temperature)

¹From Cohen et. al, "On Estimating Optimal Performance of CPU Dynamic Thermal Management." Computer Architecture Letters, Volume 2, Oct. 2003

Summary: Observations

- **Processors are becoming thermally limited**
 - It is 1st order problem ... and will just get worse.
- **Thermal solutions improve –**
But cost and form factor limit their benefit
- **Variations prohibit worst case design**
 - Too inefficient
 - Error-prone
- **Necessitates dynamic thermal management**
- **DTM evolution – from simple-reactive to smart-proactive**

Summary: Challenges

- **Reduce Power/Temperature**
 - To actively eliminate/relax heating
- **Improve measurements**
 - To know better where we are
- **Improve control**
 - To get the maximum in a given environment
- **Architecture breakthrough**
 - Significant improvement in power-awareness

The End

Backup Slides

General trends

Process Technology:

- 😊 Smaller (2x), faster (1.35x) transistors, consuming less energy per switch (1.75x) – but:
- 😞 More transistors per mm², More switches per cycle → More power density (~1.5X)

Micro-architecture:

- 😞 Traditional - More transistor switches per inst.
- 😊 New world - More efficiency → More activity per area
 - ➔ Higher power, higher power density
 - ➔ Higher temperature

The Thermal Wall

With Naïve Process Technology:

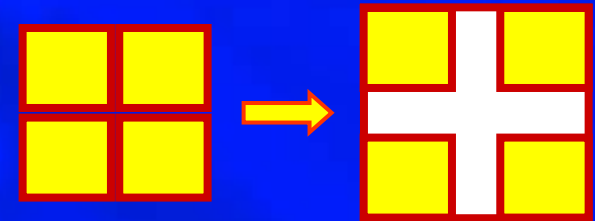
- ☺ Same Area
 - 2X transistors, 1.35X frequency
 - 1.8X-2.3X perf., 1.3X-1.5X Power, 1X Area
- ☹ Same Power →
 - 1.33X transistors, 1.35 frequency
 - 1.5X-1.8X perf., 1X Power, 0.67X Area
- ☺ Same Thermal →
 - 1.33X transistors, 0.8X power
 - 1.25X frequency, 1.35X-1.65X perf...

Simple Solutions

- Spread the transistors over the area

- ☺ Reduce power density

- ☹ Higher wire delays



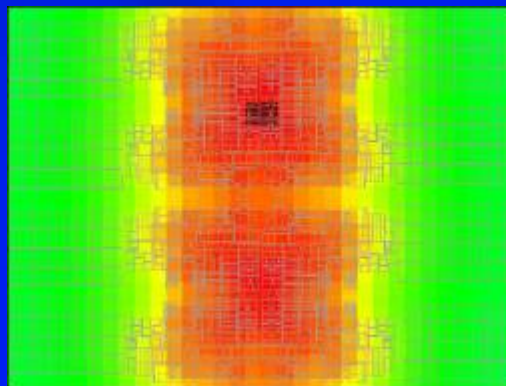
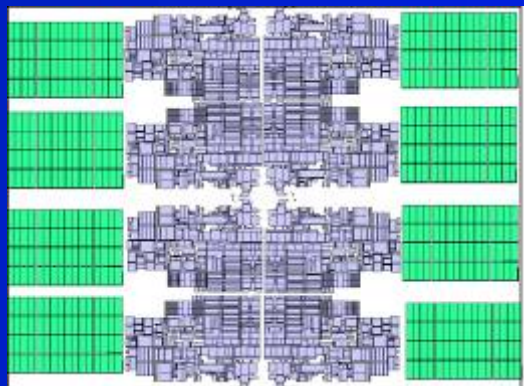
- Use voltage scaling

- ☺ Reduce power & power density

- ☺ Allow more transistors

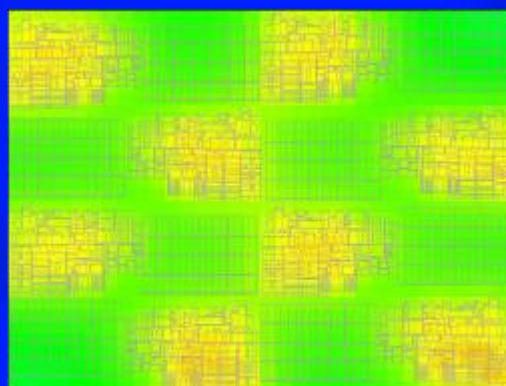
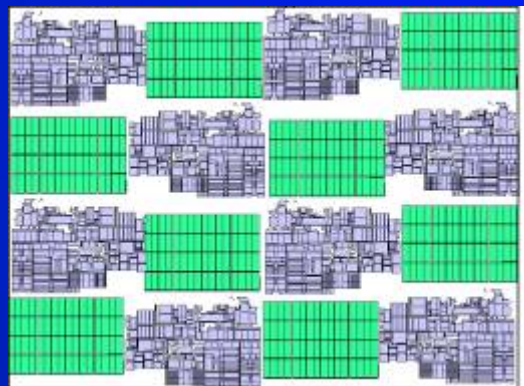
- ☹ Lower frequency → Lower performance

Thermal Variation – Space / CMP



8 cores on the center
Tjmax: 100.4°C, leakage: 19.3W

Qualitative Data



8 cores checkers layout
Tjmax: 94.0°C, leakage: 18.06W

- Uniform power distribution is better
- Chip Multi-Processor (CMP) can exploit that

r22

Replace w/ Banias
ronen; 27/4/2005

Thermal Basics

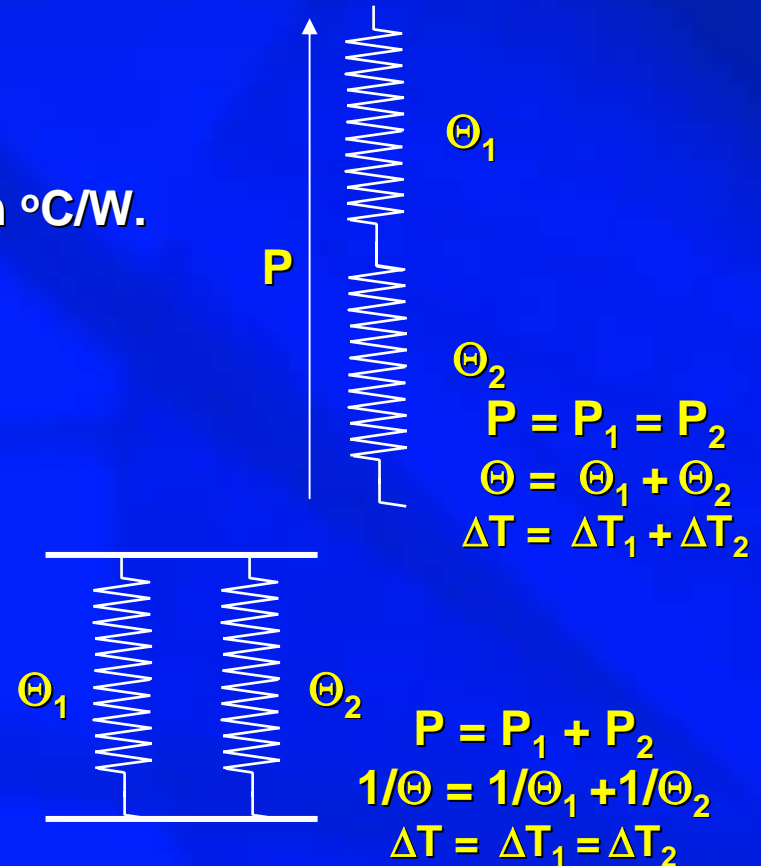
- $\Delta T = P * \Theta$
 - Θ is “Thermal resistance”
 - Similar to $\Delta V = I * R$
 - T is measured in °C, P in Watts, Θ in °C/W.

- As in resistors we see:
 - Serial: $\Theta = \Theta_1 + \Theta_2$
 - Parallel: $1/\Theta = 1/\Theta_1 + 1/\Theta_2$

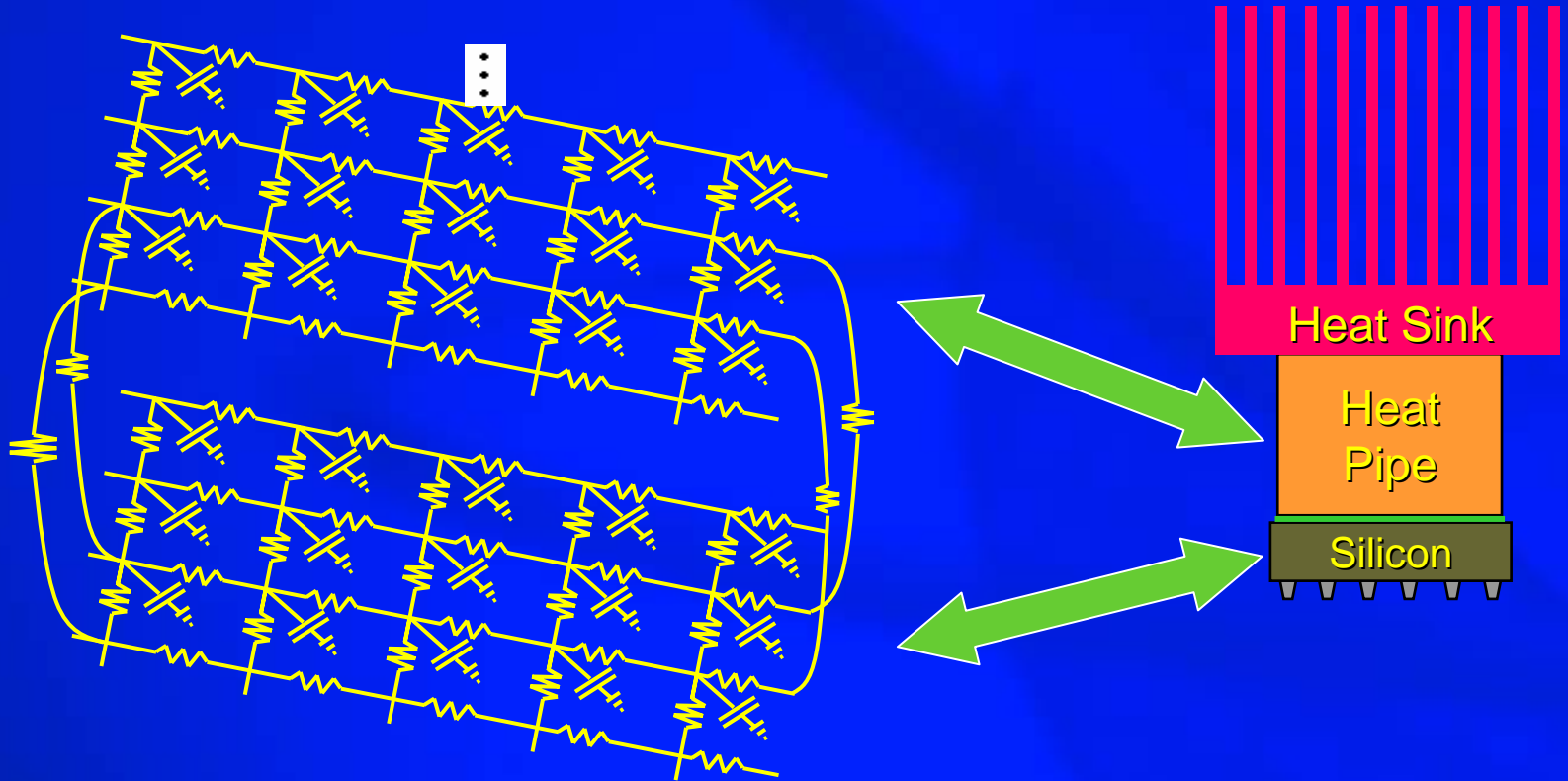
- Θ behaves like resistance
 - Grows with length
 - Reduced with area
 - For a piece of silicon, we can measure $R = \Theta$ per area: °C/(W*cm²)

→ Simplified model

- Realistic models take into account power distribution and thermal maps



RC networks for Thermal



Challenges – Improve Measurements

Simulation

- **Better, more accurate, thermal estimate**
 - More detailed and accurate power modeling
 - Integrated active/leakage with thermal modeling
 - Accurate thermal solutions

On chip

- **Temperature measurement on die.**
 - More accurate, more points on chip
 - Targeting hot spots

Challenges – Reduce Temperature

Power Awareness

- ***“Less is More”***
Do the same “work” with less power
 - Less instruction per task
 - Less micro-operation per instruction
 - Less transistor switches per micro-operation
 - Less energy per transistor switch
- **Better power/performance trade off**
 - Optimal/adaptive structure sizes, ...

Reduce power density

- **Identify hot spots – distribute them as possible**

Better cooling solutions

Challenges – Improve Control

Adaptive Thermal management

- **Eliminate risk of thermal run-away**
 - Already being done in today's processor
- **Maximize performance in thermally limited environment**
 - Technology in infancy
 - Lack of determinism is a concern for OEMs*
- **This is an optimization problem!**
 - How to control power → thermal for best performance
 - Current Challenges: SMT/MP.

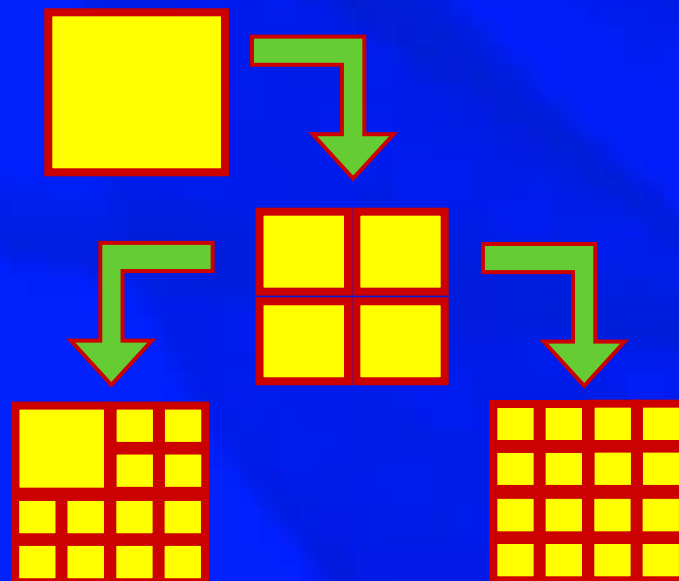
* See http://www.vanshardware.com/articles/2004/05/040517_efficeonFreeze/040517_efficeonFreeze.htm

What next? Asymmetric Cores?

- Mix big and small cores
 - Big core(s) – for single thread performance
 - Small cores – for efficient multithreaded performance*

	ST perf	MT perf
1 core	1X	1X
4 cores	$\frac{1}{2}X$	2X
16 cores	$\frac{1}{4}X$	4X
13 cores	$\frac{1}{2}X$	3.5X

Best of all worlds?



* Assume smaller core has $\frac{1}{4}$ area, $\frac{1}{4}$ power and $\frac{1}{2}$ performance